

Specification of Tokens

Lecture 2
Section 3.3

Robb T. Koether

Hampden-Sydney College

Fri, Jan 16, 2015

- 1 Alphabets and Languages
- 2 Operations on Languages
- 3 Regular Expressions
- 4 Extensions of Regular Languages
- 5 Assignment

Outline

- 1 Alphabets and Languages
- 2 Operations on Languages
- 3 Regular Expressions
- 4 Extensions of Regular Languages
- 5 Assignment

Alphabets and Strings

Definition (Alphabet)

An **alphabet** is a finite set of symbols. Traditionally, we denote an alphabet by the letter Σ .

Definition (String)

A **string** is a finite sequence of symbols.

Definition (Empty String)

The **empty string**, denoted ε , is the string that contains no symbols. The empty string has length 0.

Example (Alphabets and Strings)

- Examples:
 - The traditional alphabet is

$$\Sigma = \{A, B, C, \dots, Z\}.$$

- The binary alphabet is $\Sigma = \{0, 1\}$.
- For C programs, the alphabet is the set of ASCII characters.

Definition (Language)

A **language** is a set of (finite) strings over a given alphabet.

- A language can be (and usually is) infinite.
- The set of all even integers over the alphabet $\Sigma = \{0, 1, \dots, 9\}$ is a language.
- The set of all C programs is a language.

Outline

- 1 Alphabets and Languages
- 2 Operations on Languages**
- 3 Regular Expressions
- 4 Extensions of Regular Languages
- 5 Assignment

Operations on Languages

Definition (Language)

Let L , L_1 , and L_2 be languages.

- Union:

$$L_1 \cup L_2 = \{x \mid x \in L_1 \text{ or } x \in L_2\}.$$

- Concatenation:

$$L_1 L_2 = \{xy \mid x \in L_1 \text{ and } y \in L_2\}.$$

- Repeated concatenation:

$$L^n = L L L \cdots L \text{ } n \text{ copies of } L.$$

- Kleene closure:

$$L^* = \{\varepsilon\} \cup L \cup L^2 \cup L^3 \cdots .$$

Example (Language)

- Let

$$L_1 = \{1, 3, 5, 7, 9\}$$

and

$$L_2 = \{0, 2, 4, 6, 8\}.$$

- Describe $L_1 \cup L_2$.
- Describe $L_1 L_2$.
- Describe $(L_1 \cup L_2)^3$.
- Describe $(L_1 \cup L_2)^* L_2$.

Outline

- 1 Alphabets and Languages
- 2 Operations on Languages
- 3 Regular Expressions**
- 4 Extensions of Regular Languages
- 5 Assignment

Regular Expressions

- A **regular expression** can be used to describe a language.
- Regular expressions may be defined in two parts.
 - The basic part.
 - The recursive part.

Regular Expressions

- The basic part:
 - ε represents the language $\{\varepsilon\}$.
 - \mathbf{a} represents the language $\{\mathbf{a}\}$ for every $a \in \Sigma$.
 - Call these languages $L(\varepsilon)$ and $L(\mathbf{a})$, respectively.

Regular Expressions

- The recursive part: Let r and s denote regular expressions.
 - $r \mid s$ represents the language $L(r) \cup L(s)$.
 - rs represents the language $L(r)L(s)$.
 - r^* represents the language $L(r)^*$.
 - (r) represents the language $L(r)$.

Regular Expressions

- In other words

- $L(r \mid s) = L(r) \cup L(s)$.
- $L(rs) = L(r)L(s)$.
- $L(r^*) = L(r)^*$.
- $L((r)) = L(r)$.

Example

Example (Identifiers)

- Identifiers in C++ can be represented by a regular expression.

$$r = A \mid B \mid \cdots \mid Z \mid a \mid b \mid \cdots \mid z$$

$$s = 0 \mid 1 \mid \cdots \mid 9$$

$$t = r(r \mid s)^*$$

Regular Expressions

Definition (Regular definition)

A **regular definition** of a regular expression is a “grammar” of the form

$$\begin{aligned}d_1 &\rightarrow r_1 \\d_2 &\rightarrow r_2 \\&\vdots \\d_n &\rightarrow r_n\end{aligned}$$

where each r_i is a regular expression over $\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$.

Example

Example (Identifiers)

- We may now describe C++ identifiers as follows.

letter \rightarrow A | B | \dots | Z | a | b | \dots | z

digit \rightarrow 0 | 1 | \dots | 9

id \rightarrow *letter* (*letter* | *digit*)*

Regular Expressions

- Note that this definition does not allow recursively defined tokens.
- In other words, d_i cannot be defined in terms of d_j , not even indirectly.

Outline

- 1 Alphabets and Languages
- 2 Operations on Languages
- 3 Regular Expressions
- 4 Extensions of Regular Languages**
- 5 Assignment

Extensions of Regular Languages

Definition

We add the following symbols to our regular expressions.

- One or more instances: $r^+ = r r^*$.
- Zero or one instance: $r? = r \mid \varepsilon$.
- Character class: $[a_1 a_2 \cdots a_n] = a_1 \mid a_2 \mid \cdots \mid a_n$.

Example (Identifiers)

- Identifiers can be described as

letter \rightarrow [A-Za-z]

digit \rightarrow [0-9]

id \rightarrow *letter* (*letter* | *digit*)*

Extensions of Regular Languages

Example (Floating-point Numbers)

- Floating-point numbers can be described as

$$\textit{digit} \rightarrow [0-9]$$
$$\textit{digits} \rightarrow \textit{digit}^+$$
$$\textit{number} \rightarrow \textit{digits} (. \textit{digits})? (\text{E} [+ -]? \textit{digits})?$$

Outline

- 1 Alphabets and Languages
- 2 Operations on Languages
- 3 Regular Expressions
- 4 Extensions of Regular Languages
- 5 Assignment**

Assignment

Assignment

- Read Section 3.3.
- Exercises 2, 3.